

This Page Is Inserted by IFW Operations
and is not a part of the Official Record

BEST AVAILABLE IMAGES

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images may include (but are not limited to):

- BLACK BORDERS
- TEXT CUT OFF AT TOP, BOTTOM OR SIDES
- FADED TEXT
- ILLEGIBLE TEXT
- SKEWED/SLANTED IMAGES
- COLORED PHOTOS
- BLACK OR VERY BLACK AND WHITE DARK PHOTOS
- GRAY SCALE DOCUMENTS

IMAGES ARE BEST AVAILABLE COPY.

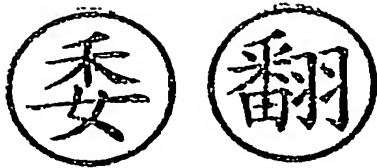
**As rescanning documents *will not* correct images,
please do not report the images to the
Image Problem Mailbox.**

[Number of appeal against examiner's decision
of rejection]

[Date of requesting appeal against examiner's
decision of rejection]

[Date of extinction of right]

Copyright (C); 1998,2003 Japan Patent Office



02RL/48

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平7-271792

(43) 公開日 平成7年(1995)10月20日

(51) Int.Cl.

G 0 6 F 17/27

識別記号

庁内整理番号

F I

技術表示箇所

8219-5L

G 0 6 F 15/ 38

E

審査請求 未請求 請求項の数4 O L (全 18 頁)

(21) 出願番号

特願平6-61527

(22) 出願日

平成6年(1994)3月30日

(71) 出願人 000004226

日本電信電話株式会社

東京都千代田区内幸町一丁目1番6号

(72) 発明者 永田 昌明

東京都千代田区内幸町1丁目1番6号 日

本電信電話株式会社内

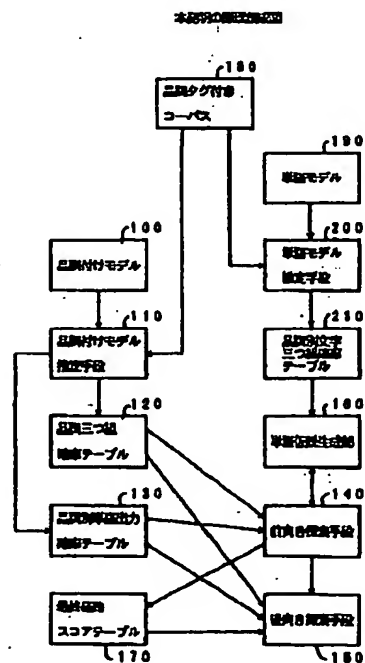
(74) 代理人 弁理士 伊東 忠彦

(54) 【発明の名称】 日本語形態素解析装置及び日本語形態素解析方法

(57) 【要約】

【目的】 本発明の目的は、単語分割と品詞付与の組を、入力文が未知語を含む場合でも最も尤もらしい候補から順に任意の数だけ出力できるので、高い精度をもち、頑強でかつ柔軟なインタフェースを持つ日本語形態素解析装置及び日本語形態素解析方法を提供することである。

【構成】 本発明は、統計的な言語モデルである品詞付けモデルに基づいて、多重マルコフ過程の縮退、前向き動的計画法、内路向きA*探索を用いて、二回探索することにより、辞書に未登録の単語の形態素解析も含めて実現する。



【特許請求の範囲】

【請求項 1】 与えられた日本語の入力文を形態素解析する装置であって、
連続する 3 つの品詞の組において直前の 2 つの品詞が与えられた時の 3 つ目の品詞の確率である品詞三つ組確率と、品詞が与えられた時の単語の確率である品詞別単語出力確率から、文を構成する単語列と各単語に付与された品詞列の同時確率を与える品詞付けモデルと、単語に分割され、かつ、品詞が付与された文の集合から、品詞三つ組確率と品詞別単語出力確率を推定する品詞付けモデル推定手段と、
該品詞三つ組確率を記憶する品詞三つ組確率テーブルと、
該品詞別単語出力確率を記憶する品詞別単語出力確率テーブルと、
文頭からある単語に至るまでの単語列と品詞列の同時確率の最大値を記憶する最適経路スコアテーブルと、
該最適経路スコアテーブルに記憶されている該文頭からある単語に至るまでの単語列と品詞列の同時確率の最大値を、該単語を含む直前の 2 つの単語に付与された品詞の組が異なる場合毎に、該単語を含む直前の 3 つの単語に付与された品詞三つ組確率、該単語の品詞別単語出力確率、及び文頭からその単語の直前の単語に至るまでの単語列と品詞列の同時確率の最大値から求め、該最適経路スコアテーブルに記録する前向き探索手段と、
単語列と品詞列の同時確率を最大化するような入力文の単語分割及び品詞付与の候補を、文末からある単語に至るまでの単語列と品詞列の同時確率の最大値、及び該最適経路スコアテーブルに記録された、文頭からその単語に至るまでの単語列と品詞列の同時確率の最大値から、最も尤もらしい候補から順番に任意の個数だけ求める後向き探索手段とを有し、
該品詞三つ組確率と該品詞別単語出力確率から構成される該品詞付けモデルに基づいて、まず、該前向き探索手段を用いて、文頭からある単語までの単語列と品詞列の同時確率の最大値を該最適経路スコアテーブルに記憶し、該最適経路スコアテーブルの値に基づいて該後向き探索手段を用いて最も尤もらしい順番に任意の個数の形態素解析候補、即ち、単語列と品詞列の組を求めることを特徴とする日本語形態素解析装置。
【請求項 2】 単語に分割され、かつ品詞が付与された文の集合から、品詞別文字三つ組確率を推定する単語モデル推定手段と、
該品詞別文字三つ組の確率から、該単語毎に品詞別出力確率を与える単語モデルと、
該品詞別文字三つ組確率を記憶する品詞別文字三つ組確率テーブルと、
ある文字位置から始まる入力文の部分文字列の中から、品詞別文字三つ組確率に基づいて、最も尤もらしい順番に任意の個数の単語仮説、即ち、表記、品詞及び品詞別

出力確率の組を求める単語仮説生成部と、
該入力文中に未知語が存在する場合にも、前記前向き探索手段において、該品詞別文字三つ組から構成される該単語モデルを用いて、文字表記から単語仮説を生成し、
前記後向き探索手段において、文全体を考慮して最も尤もらしい、未知語の単語区切り、品詞及び品詞別単語出力確率を求めることを含む請求項 1 記載の日本語形態素解析装置。

【請求項 3】 日本語の文が入力されると、連続する 3 つの連続する品詞の組において直前の二つの品詞が与えられた時の三つ目の品詞の確率である品詞三つ組確率と品詞が与えられた時の単語の確率である品詞別単語出力確率に基づいて、文頭から文末へ一文字ずつ進む動的計画法を用いて、文頭からある単語に至るまでの単語列と品詞列の同時確率の最大値を求める前向き探索を行い、
文末から文頭へ進む A* アルゴリズムを用いて、文を構成する単語列と各単語に付与された品詞列の同時確率を最大化する形態素解析候補、即ち、単語列と品詞列の組を最も尤もらしい順に一つずつ求める後向き探索を行う日本語形態素解析方法。

【請求項 4】 入力文が辞書に登録されていない未知語を含む場合に、前記前向き探索において、前記品詞別文字三つ組確率に基づいて、単語を構成する文字列の始まりと終わり、品詞、品詞別出力確率からなる単語仮説を生成し、
前記後向き探索により最も尤もらしい単語分割と品詞付与を決定する請求項 3 記載の日本語形態素解析方法。

【発明の詳細な説明】

【0001】

【産業上の利用分野】 本発明は、日本語形態素解析装置及び日本語形態素解析方法に係り、特に、日本語文を単語に分割し、各単語に品詞を付与する日本語形態素解析装置及び日本語形態素解析方法に関する。

【0002】 詳しくは、品詞タグ付きコーパスから統計的な手法によって求めた言語モデルを用いることにより、文法や辞書を対象領域に自動的に適応化し、入力文の長さに比例する計算量で効率的に精度よく形態素解析を行い、任意の個数の最尤な形態素解析候補を求め、入力文が未知語を含む場合にも形態素解析が行えるようにしたものである。

【0003】

【従来の技術】 従来の日本語形態素解析システムは、文法規則と発見的探索に基づく方式が主流であり、これらは、文法規則（言語モデル）として品詞接続表を用いて、形態素解析候補の探索の際には、発見的規則（ヒューリスティックス (heuristics)）を用いて候補の順位付けを行う（第 1 の従来の方法）。

【0004】 形態素解析で用いられる発見的規則としては、入力文と照合する最も長い辞書中の単語を含む形態素解析候補を優先する最良一致法（第 1 の従来の方法）

や、文節数が最も少ない形態素解析候補を優先する文節数最小法（吉村・日高・吉田「文節数最小法を用いたべきき日本語文の形態素解析」情処論 Vol. 24 No. 1, pp. 40-46, 1983）がある（第2の従来の方法）。

【0005】また、文節数最小法を拡張し、解を詳細に順序付けする方法として、接続コスト最小法（久光・新田「接続コスト最小による形態素解析の提案と計算量の評価について」信学会 NLC90-8, pp. 17-24, 1990）がある。これは、接続コストが最小となる形態素解析候補を、動的計画法を用いて、入力文の長さに比例する計算量を求めることができるアルゴリズムである（第3の従来の方法）。

【0006】また、近年、上記の従来の方法である「規則に基づく方法」に対して、「統計に基づく方法」が研究され始めている。これは、必ずしも根拠が明確ではなかった発見的な接続コストの代わりに、大量のテキストデータから求めた統計的言語モデルに基づく接続確率を使用する方法である（第4の従来の方法）。

【0007】また、上記の第3の実施例の接続コスト最小法における接続コストの代わりに、確率正規文法に基づくコストを使用する方法が提案されている（松延・日高・吉田「確率文節文法による構文解析」情処 NL56-3, 1986）。この方法では、確率正規文法、即ち、品詞の接続確率と品詞別の単語の出力確率の情報を接続コストの計算に用いる。

【0008】上記の統計的な言語モデルを用いることにより、形態素解析候補の優先度の根拠が明確により、かつ、候補の詳細な順序付けができるようになる。

【0009】また、従来の日本語形態素解析システムにおいて辞書に登録されていない未知語の扱いは、非常に「場当たりの（ad hoc）」である。多くのシステムは、「同じ字種の文字の連続が単語を構成することが多い」、或いは、「カタカナの連続は外来語の名詞であることが多い」というような字種に関する発見的規則を用いて、単語の分割及び、品詞の付与を行っている（吉村・武内・津田・首藤「未登録語を含む日本語文の形態素解析」情処論 Vol. 30 No. 3, pp. 294-301, 1989）。或いは、付属語列等から文節を推定し、そこから、付属語を取り除いた部分列を未知語とみなす方法が使われていることが多い（第5の従来の方法）。

【0010】未知語問題に対する統計的な解決策として、造語単位の2つの組を用いて単語の生起確率を推定する方法が提案されている（永井・日高「日本語における単語の造語モデルとその評価」情処論 Vol. 34 No. 9, pp. 1944-1955, 1993）。この方法は、単語辞書の見出し語から造語モデルのパラメータを推定する。

【0011】

【発明が解決しようとする課題】しかしながら、上記従来の形態素解析方法には、以下のような問題点がある。

【0012】上記の第1の従来の方法の最長一致法に

は、

（1）単語または文節を認識する段階で複数の候補が存在する場合に、入力字列の長さが長いものに高い優先順位を与えることの根拠が明らかではない；

05 （2）単語または文節の長さという局所的な評価しているために、文全体に対する尤度の評価ができない；という問題がある。

【0013】さらに、上記の第2の従来の方法である文節数最小法は、

10 （1）入力文字列を構成する文節の総数が最小になる解釈を優先する根拠が明らかでない；

（2）一般に、文節数最小となる候補は複数存在するので、さらに何らかの手段で候補を詳細に順位付けする必要がある；等の問題がある。

15 【0014】さらに、上記第3の従来の方法である接続コスト最小法は、具体的な接続コストの設定方法には言及しておらず、「自立語と付属語の接続コストは小さい」というような内容を表す接続コストが発見的な方法で与えられている。

20 【0015】また、従来提案された確率正規文法に基づく日本語の形態素解析法には、次のような問題がある。

【0016】（1）一般に、品詞の接続確率（品詞の二つ組の確率）だけでは、言語の局所的な性質を十分にモデル化できないことが多い。そこで、英語の確率的形態素解析では、品詞の三つ組の確率を用いることにより、精度の向上を達成している。しかし、接続コスト最小法及び確率文節文法、単語または品詞の接続（即ち、二つ組）に対してコスト与えられることを想定しており、三つ組の確率に基づくコストは扱えない。

30 【0017】（2）接続コスト最小法は、最尤候補を1つだけ出力する。出力が上位N個の最尤候補を含むように、接続コスト最小法を拡張することは可能であるが、その場合には、予め、候補Nを決めておく必要がある。従って、例えば、形態素解析の出力を構文解析の入力として利用する場合などに、状況に応じて、最も尤もらしい候補から順番に任意の個数の候補を取り出すようなことはできない。

40 【0018】また、上記の未知語の扱いの方法における前者の方法は、文字列の単語らしさを評価する方法であり、上記の後者の方法は、ある文脈における文字列の単語らしさを評価する方法とみなすことができるが、どちらの場合も尤度の根拠が不明確であり、単語仮説の詳細な順位付けも難しい。

【0019】上記第6の従来の方法は、未知語の単語らしさを評価する有力な手段を与えるが、次のような問題点がある。

【0020】（1）漢字で表記される複合語に対する造語モデルであり、それ以外では、造語単位の設定基準が必ずしも明確ではなく、自動的な処理が難しい。

50 【0021】（2）生起確率の推定が目標であり、品詞

を考慮していない。従って、品詞の推定及び、品詞別の出力確率の推定ができない。

【0022】(3) 辞書の見出し語からモデルパラメータを推定するので、モデルパラメータが対象領域のテキストの性質を反映しない。

【0023】(4) 文字列の単語らしさを評価するだけで、その文字列が現れた文脈(前後の文字列)の情報を使って、単語としても尤もらしさを評価しているわけではない。

【0024】本発明は、上記の点に鑑みなされたもので、上記従来の問題点を解決し、単語分割と品詞付与の組を、入力文が未知語を含む場合でも最も尤もらしい候補から順に任意の数だけ出力できるので、高い精度をもち、頑強でかつ柔軟なインタフェースを持つ日本語形態素解析装置及び日本語形態素解析方法を提供することを目的とする。

【0025】

【課題を解決するための手段】図1は、本発明の原理構成図である。

【0026】本発明の形態素解析装置は、日本語の入力文が与えられ、形態素解析を行う装置であって、直前の2つの品詞が与えられた時の3つ目の品詞の確率である品詞三つ組確率と、品詞が与えられた時の単語の確率である品詞別単語出力確率から、文を構成する単語列と各単語に付与された品詞列の同時確率を与える品詞付けモデル100と、単語に分割され、かつ、品詞が付与された文の集合から、品詞三つ組確率と品詞別単語出力確率を推定する品詞付けモデル推定手段110と、品詞三つ組確率を記憶する品詞三つ組確率テーブル120と、品詞別単語出力確率を記憶する品詞別単語出力確率テーブル130と、文頭からある単語に至るまでの単語列と品詞列の同時確率の最大値を記憶する最適経路スコアテーブル170と、最適経路スコアテーブル170に記憶されている文頭からある単語に至るまでの単語列と品詞列の同時確率の最大値を、単語を含む直前の2つの単語に付与された品詞の組が異なる場合毎に、単語を含む直前の3つの単語に付与された品詞三つ組確率、単語の品詞別単語出力確率、及び文頭からその単語の直前の単語に至るまでの単語列と品詞列の同時確率の最大値から求め、最適経路スコアテーブル170に記録する前向き探索手段140と、単語列と品詞列の同時確率を最大化するような入力文の単語分割及び品詞付与候補を、文末からある単語に至るまでの単語列と品詞列の同時確率の最大値、及び最適経路スコアテーブル170に記録された、文頭からその単語に至るまでの単語列と品詞列の同時確率の最大値から最も尤もらしい候補から順番に任意の個数だけ求める後向き探索手段150とを有し、品詞三つ組確率と品詞別単語出力確率から構成される品詞付けモデルに基づいて、まず、前向き探索手段140を用いて、文頭からある単語までの単語列と品詞列の同時確

率の最大値を最適経路スコアテーブル170に記憶し、最適経路スコアテーブル170の値に基づいて後向き探索手段150を用いて最も尤もらしい順番に任意の個数の形態素解析候補を求める。

05 【0027】また、本発明の日本語解析装置は、単語に分割され、かつ品詞が付与された文の集合から、品詞別文字三つ組の確率を推定する単語モデル推定手段200と、品詞別文字三つ組の確率から、単語の品詞別出力確率を与える単語モデル190と、品詞別文字三つ組確率を記憶する品詞別文字三つ組確率テーブル210と、ある文字位置から始まる入力文の部分文字列の中から、品詞別文字三つ組確率に基づいて、最も尤もらしい順番に任意の個数の単語仮説、即ち、表記、品詞及び品詞別出力確率の組を求める単語仮説生成部160と、入力文中
10 に未知語が存在する場合にも品詞別文字三つ組から構成される単語モデルを用いて、文字表記から単語仮説を生成し、後向き探索手段150において、文全体を考慮して最も尤もらしいように未知語の単語区切り、品詞及び品詞別単語出力確率を求める。

20 【0028】図2は、本発明の原理を説明するためのフローチャートである。

【0029】本発明の日本語形態素解析方法は、品詞三つ組確率と品詞別単語出力確率から求められる、文を構成する単語列と単語に付与された品詞列の同時確率を最大とするような形態素解析候補を、文頭から文末に探索する前向き探索と文末から文頭に探索する後向き探索を用いて、最も尤もらしい順番に任意の個数だけ求める。
25 また、上記の前向き探索は動的計画法を用いて、文頭からある単語に至るまでの単語列と品詞列の同時確率の最大値を求め、記録する。

30 【0030】また、上記の後向き探索はA*アルゴリズムを用いて、文末からある単語に至るまでの単語列と品詞列の同時確率の最大値と、前向き探索で記録された、文頭からある単語に至るまでの単語列と品詞列の同時確率の最大値により形態素解析候補の順位付けを行う。

35 【0031】さらに、本発明は、入力文が辞書に登録されていない未知語を含む場合に、前向き探索において、単語を構成する文字列の始まりと終わり、品詞及び品詞別単語出力確率からなる単語仮説を、品詞別文字三つ組確率に基づいて入力文字列から生成し、後向き探索により、最も尤もらしい単語分割と品詞付与を決定することにより形態素解析候補を求める。

【0032】

45 【作用】本発明は、シンボルの二次の拡大を考えて二重マルコフ過程を一重マルコフ過程に縮退させることにより、品詞三つ組確率を、従来の接続コスト最小法を拡張した枠組の中で扱えるようにする。具体的には、文頭から文末へ文字ずつ進む動的計画法を用いて、単語列と品詞列の同時確率、即ち、品詞三つ組確率と品詞別単語出力確率の積を最大化するような形態素解析候補(単語
50

分割と品詞付与の組)を求める。この探索は、「前向き探索」と呼ばれる。前向き探索では、文頭からある文字位置までの入力文の部分列に対する単語系列と品詞系列の同時確率の最大値が、最適経路スコアテーブルに、最後の単語と最後の二重マルコフ過程の状態の異なる組み合わせ毎に記録される。

【0033】次に、前向き探索で作成した最適経路スコアテーブルに基づいて、文末から文頭へ進むA*探索を用いて、最も尤もらしい順に一つずつ形態素解析候補を求める。この探索は、「後向き探索」と呼ばれる。後向き探索では、文末からある単語に至るまでの部分経路(単語列と品詞列の組)の候補が完全な経路のコスト、即ち、後向きの部分経路のコストと残りの最適部分経路のコストの和、に基づいて順位付けされる。残りの最適部分経路、即ち、その単語から文頭に至るまでの最適部分経路のコストは、前向き探索において、最適経路スコアテーブルに記録されているので、A*探索の性質により、この方法で必ず最適解が求められることが保証されている。また、最適解が求められた後に、この解を取り除き、さらに、A*探索を続けることにより、次の最適解が求められることも保証されている。

【0034】従って、このような、マルコフ過程の縮退、前向きの動的計画法、及び後向きのA*探索の組み合わせにより、(1)品詞の三つ組の確率を用いて、(2)最も尤もらしい順に任意の個数の候補を求める、日本語形態素解析装置が実現できる。また、この方法は、任意の高次のマルコフ過程を扱うように、自然に拡張できる。

【0035】未知語を推定する問題は、(1)単語境界(単語を構成する文字列の始まりと終わり)、(2)品詞別の出力確率、からなる三つ組を求める問題に帰着される。この三つ組のことを「単語仮説」と呼ぶことにする。入力文字列から辞書に収録されていない単語仮説を生成する手段があれば、入力文が未知語を含む場合でも、文全体を考慮して最も尤もらしい形態素解析候補を求めることができる。

【0036】単語を構成する文字列と品詞が与えられた時に、その単語の品詞別の出力確率を計算するモデルのことを「単語モデル」と呼ぶことにする。単語モデルがあれば、入力文の任意の部分文字列の単語らしさを判定することができる。

【0037】本発明では、品詞別文字三つ組確率から構成される単語モデルを用いて、入力文に対して単語仮説を生成する。なお、単語モデルのパラメータは、品詞タグ付きコーパスから推定する。具体的には、前向き探索では、入力文の全ての文字位置において、その最左部分文字列と一致し、かつ、辞書に登録された単語を提案する。これに加えて、未知語を扱うために、最左部分文字列と一致するが、辞書には登録されていない単語仮説を、予め決めた上限以下の個数、最も尤もらしいものか

ら順に、提案する。次に、後向き探索では、文全体を考慮に入れて最も尤もらしい単語分割と品詞付与を決定する。

【0038】従って、このような、統計的な単語モデルに基づく文字表記からの単語仮説生成法、及び、上記前向き探索と後向き探索から構成される日本語形態素解析法の組み合わせにより、(1)未知語の単語らしさを自動的に評価でき、(2)品詞推定が可能で、(3)対象とするテキストデータの性質を反映し、(4)文全体から見ても最も妥当な解釈をする、日本語形態素解析装置が実現できる。

【0039】

【実施例】以下、図面と共に本発明の実施例を詳細に説明する。

【0040】図3は、本発明の一実施例の概略ブロック図を示す。同図に示すように、本発明の形態素解析装置は、前向きDP探索部1、後向きA*探索部2、品詞タグ付きコーパス3、品詞付けモデル推定部4、品詞三つ組確率テーブル5、品詞別単語出力確率テーブル6、単語モデル推定部7、品詞別文字三つ組確率テーブル8、単語仮説生成部9及び、最適経路スコアテーブル10より構成される。

【0041】前向きDP探索部1は、入力文の文頭から文末へ文字ずつ進む動的計画法(DP:Dynamic Programming)を用いて、単語列と品詞列の同時確率、即ち、品詞三つ組確率と品詞別単語出力確率の積を最大化するような、入力文の単語分割と品詞付与の組を求める。品詞二つ組確率テーブル5及び品詞別単語出力確率テーブル6は、それぞれ、品詞三つ組確率及び品詞別単語出力確率を格納する。この前向きの動的計画法では、文頭からある単語に至るまでの単語列と品詞列の同時確率を最大とする最適部分経路(単語列と品詞列の組)の確率が、最後の単語と最後の二重マルコフ過程の状態の異なる組み合わせ毎に計算され、最適経路スコアテーブル10に記録される。

【0042】単語仮説生成部9は、前向きの動的計画法において、入力文のある文字位置における最左部分列と照合する単語を検索する際に、最左部分列を単語表記とする単語仮説を最も尤もらしい順に、予め定めた個数以下だけ生成する。

【0043】後向きA*探索部2は、前向きDP探索部1で計算された最適経路スコアテーブル10の内容を入力とし、文末から文頭へ単語ずつ進むA*アルゴリズムを用いて、品詞三つ組確率と品詞別単語出力確率の積が最も大きいものから順番に一つずつ形態素解析候補を求める。

【0044】品詞付けモデル推定部4は、品詞タグ付きコーパス3、即ち、単語に分割され、かつ、品詞が付与された大量の日本語文のデータから、品詞三つ組確率と品詞別単語出力確率を推定し、品詞三つ組確率テーブル

5 及び品詞別単語出力確率テーブル 6 へ格納する。

【0045】単語モデル推定装置 7 は、品詞タグ付きコーパスから品詞別文字三つ組確率を推定し、品詞別文字三つ組確率テーブル 8 へ格納する。

【0046】以下では、品詞付けモデルの推定、マルコフ過程の縮退、動的計画法を用いた前向き探索、A* アルゴリズムを用いた後向き探索、単語モデルの推定、単語仮説の生成の順に説明する。

【0047】(1) 品詞付けモデルの推定

品詞付けモデル推定部 4 が品詞タグ付きコーパスから、品詞三つ組確率と品詞別単語出力確率から構成される品詞付けモデルを推定する手順を示す。

【0048】本発明では、日本語の形態素解析のため

$$P(W, T) = \prod_{i=1}^n P(t_i | t_{i-2}, t_{i-1}) P(w_i | t_i)$$

(1)

【0051】実用的には、文の境界が重要な情報を担っているため、式(1)の代わりに式(2)を用いるのが有効である。ここで“#”は文頭及び文末を表す特別な記号

$$P(W, T) = P(t_1 | \#) P(w_1 | t_1) P(t_2 | \#, t_1) P(w_2 | t_2)$$

$$\prod_{i=1}^n P(t_i | t_{i-2}, t_{i-1}) P(w_i | t_i) P(\# | t_{n-1}, t_n)$$

(2)

【0053】品詞三つ組確率 $P(t_i | t_{i-2}, t_{i-1})$ と品詞別単語出力確率 $P(w_i | t_i)$ は、次式(3)により品詞タグ付きコーパスにおける相対頻度から推定できる。但し、 f は相対頻度を求める関数、 $N(w, t)$ は品詞タグ t を持つ単語 w が現れた回数、 $N(t_{i-2}, t_{i-1}, t_i)$ は、品詞タグ列 t_{i-2}, t_{i-1}, t_i が現れた回数である。

$$P(t_i | t_{i-2}, t_{i-1}) = f(t_i | t_{i-2}, t_{i-1})$$

$$= \frac{N(t_{i-2}, t_{i-1}, t_i)}{N(t_{i-2}, t_{i-1})}$$

(3)

$$P(w_i | t_i) = f(w_i | t_i) = \frac{N(w, t)}{N(t)}$$

(4)

【0055】一般に、相対頻度から、直接、品詞三つ組確率を推定するのに十分な量の品詞タグ付きコーパスを用意するのは難しい。そこで、品詞三つ組確率の平滑化のために、内挿推定法を用いる。これは、式(5)に示す

$$P(t_i | t_{i-2}, t_{i-1}) = q_3 f(t_i | t_{i-2}, t_{i-1}) + q_2 f(t_i | t_{i-1}) + q_1 f(t_i) + q_0 V \quad (5)$$

ここで、 f は相対頻度を求める関数であり、 V はすべての品詞タグが一様に出現する確率である。非負の重み計

に、品詞三つ組モデル(トライポスモデル(tri-POS))と呼ばれる統計的言語モデルを用いる。

【0049】まず、入力文が単語系列 $W = w_1, w_2, \dots, w_n$ に分割され、各単語に品詞系列 $T = t_1, t_2, \dots, t_n$ が付与されたとする。トライポス(以下tri-POS)

05 モデルでは、式(1)に示すように、単語列と品詞列の同時確率 $P(W, T)$ を、品詞三つ組確率 $P(t_i | t_{i-2}, t_{i-1})$ と品詞別単語出力 $P(w_i | t_i)$ の積で近似する。従って、日本語の形態素解析は、品詞列と単語列の同時確率を最大化する単語分割と品詞付与の組を見つける問題に帰着される。

【0050】

【数1】

である。

【0052】

【数2】

【0054】

【数3】

ように、一つ組、二つ組、三つ組の相対頻度を内挿することにより、三つ組の確率を求める方法である。

【0056】

数 q_i は $q_3 + q_2 + q_1 + q_0 = 1$ を満足する。内挿

50 推定法では、観測データの確率が最大となるように、重

み計数が決定される。

【0057】図4は、品詞タグ付きコーパスの一例を示す。この例では、各文に対して文識別番号が与えられ、各文は、空白文字により単語に区切られている。各単語は、その表記に続いて品詞が付与されている。ここで、“/”は表記の品詞の区切りを示す記号である。例えば、『もしもし通訳電話国際会議事務局ですか?』については、

“3001-100

もしもし/ 感動詞、/ 記号 通訳電話国際会議事務局/ 固有名詞 です/ 助動詞・終止 か/ 終助詞、?/記号”となる。

【0058】図5は、品詞三つ組確率の一例である。品詞三つ組確率のデータは三つの品詞とその確率の4つの

$$P(u_i | u_{i-1}) = P(t_i | t_{i-2}, t_{i-1}) \quad (6)$$

式(6)を式(1)に代入すると、式(7)が得られる。

【0063】

$$P(W, T) = \prod_{i=1}^n P(u_i | u_{i-1}) P(w_i | t_i) \quad (8)$$

【0064】式(7)は、通常の一次マルコフ過程と同じ形式である。このように、結合状態を考えることにより、任意の高次マルコフ過程は、一次マルコフ過程へ還

$$P(W_i, T_i) = P(W_{i-1}, T_{i-1}) P(u_i | u_{i-1}) P(w_i | t_i)$$

(8)

この式(8)から、各結合状態 u_i に対して単語列と品詞列の同時確率 $P(W_i, T_i)$ の最大値を求めるには、(1)各結合状態 u_{i-1} に対して単語列と品詞列の同時確率 $P(W_{i-1}, T_{i-1})$ の最大値を記憶し、(2)各結合状態 u_{i-1} に対する同時確率 $P(W_{i-1}, T_{i-1})$ の最大値と式(8)を用いて全ての結合状態系列 u_i と全ての部分単語系列 w_i について同時確率 $P(W_i, T_i)$ を計算し、各結合状態系列 u_i に対する同時確率 $P(W_i, T_i)$ の最大値を選択すればよいことがわかる。

【0065】従って、 i を1から n まで1ずつ増やすことにより、同時確率 $P(W_n, T_n)$ を最大化する結合状態系列 u_n を求めることができる。しかし、日本語は単語を分かち書きする習慣がないので、入力文の単語分割を事前に一通りには決めることができない。そこで、前向き探索は、動的計画法を用いて、様々な単語分割の可能性を考慮しながら、同時確率 $P(W_n, T_n)$ を最大化する結合状態系列 u_n を求められるように工夫されている。

【0066】次に、同時確率 $P(W_n, T_n)$ を最大化する結合状態系列 u_n が求めれば通常のビタビアルゴリズムと同様に、この結合状態系列 u_n へ至るまでの経路を逆に辿ることにより、同時確率 $P(W_n, T_n)$ を最大化する結合状態系列 $U = u_1, u_2, \dots, u_n$ と単語系列

要素から構成されるリスト構造で表されている。

【0059】図6は、品詞別単語出力確率の一例である。品詞別単語出力確率のデータは、品詞、表記、確率の3つの要素から構成されるリスト構造で表される。

05 【0060】(2) マルコフの過程の縮退

以下では、前向きDP探索と後向きA*探索を用いてN-Bestの形態素解析候補を求めるアルゴリズムが、式(1)で表される品詞モデルに対して、二次マルコフ過程を一次マルコフ過程に縮退させる操作を施すことにより導かれることを示す。

【0061】 $u_i = t_i$ 及び $u_i = t_{i-1} t_i$ とすることにより、結合状態系列 $U = u_1, u_2, \dots, u_n$ を定義すると、次の関係が成り立つ。

【0062】

【数4】

元することができる。部分単語系列を $W_i = w_1, w_2, \dots, w_n$ 、部分品詞系列を $T_i = t_1 \dots t_i$ と定義すると、次のような関係がある。

30 $W = w_1 \dots w_n$ が得られる。この結合状態系列において、各結合状態の最初の状態を無視すれば、同時確率 $P(W_n, T_n)$ を最大化する状態系列(品詞タグ系列) $T = t_1 \dots t_n$ が得られる。これが、後向き探索である。後向き探索では、前向き探索で求めた部分形態素列の確率の最大値をA*アルゴリズムのヒューリスティック関数として用いることにより、最も尤もらしい順に一つずつ形態素解析候補を取り出すことができる。

35 【0067】(3) 動的計画法を用いた前向き探索

前向きDP探索部1が最適経路スコアテーブル10に格納する値を計算する手順を示す。最初に前向きDP探索部で用いられるデータ構造及び、補助的な関数について説明し、続いて具体的な処理の流れを説明する。

40 【0068】まず、図7に示すようなスロット(フィールド)を持つ、“parse”と“word”という2つの構造体(レコード)を定義する。構造体parseは、最適経過スコアテーブルにおいて、単語の情報及び文頭からその単語へ至る最適部分経路(同時確率が最大となるような部分単語列と部分品詞列の組、即ち形態素列)の情報を格納するのに用いられる。パーススタート“parse. start”とパースエンド“parse. end”は、入力文における単語の開始位置と、終了位置のインデックスである。

45 “parse. pos”は、単語の品詞で、ここでは、単語の品詞、活用型、活用形のリストを用いている。“Parse. nt

h-order-state”は、この単語を含む最後の2つの単語の品詞のリストである。このスロットは、二重マルコフ過程の結合状態に対応する。“Parse. prob-so-far”は、文頭から現在の単語に至るまでの最適部分経路のスコアである。構造体wordは、辞書(品詞別単語出力確率テーブル6)において、個々の単語の情報を格納するのに用いられる。“Word. form, word. pos, word. prob”は、それぞれ、単語の表記、品詞、品詞別の出力確率を表す。

【0069】最適経路スコアテーブル(以下ではパーステーブルと呼ぶ)は、最後の単語と最後の二重マルコフ過程の状態、即ち、“parse. start, parse. end, parse. nth-order-state”の値のリストをキーとし、同じキーを持つパース(parse)構造の中で、最適部分経路スコアが最良なもの(複数ある場合はそのリスト)を値として保持するテーブルである。

【0070】このパーステーブルの操作関数として、レジスタパース(resister-parse)とゲットパース(get-parses)を定義する。関数レジスタパース(resister-parse)はパース構造と最適経路スコアテーブルを引数として、そのパース構造をパーステーブルに登録する。この場合、同じキーを持つパース構造が既に登録されている場合は、最適部分経路のスコアが同じならば、両方ともテーブルに残す。関数ゲットパース(get-parses)は、入力文の文字位置と最適経路スコアテーブルを引数とし、その文字位置が単語の終了位置となっている。パーステーブルの中の全ての要素(パース構造)を集め、最後の二つの品詞(parse. nth-order-state)が同じパース構造の中で、最良の最適部分の経路のスコア(parse. prob-so-far)を持つパース構造のみを集めたのを返す。

【0071】図8は、本発明の一実施例の前向きDP探索部の動作を説明するためのフローチャートである。以下では、図7に従って前向きDP探索部1の動作を説明する。まず、パーステーブルには、文頭に対応する特別なパース構造を予め登録し、入力文字列の最後には、文末に対応する特別な記号を付加する。

【0072】前向きDP探索は、入力文の先頭から始まり、文末方向へ文字ずつ進む。

【0073】ステップ1) 入力文の文字位置を表す変数iを0に設定する。

【0074】i=0

ステップ2) 探索が文末に達したかを判断する。ここで、関数lengthは、文字列の長さを返す関数である。もし、文末に達していれば、前向きDP探索を終了する。そうでなければ以下の処理を各文字位置で行う。

【0075】i < length(string)

ステップ3) まず、関数get-parseを用いて単語の終了位置が文字位置iであるパース構造の中で、最後の二つの単語の品詞の異なる組み合わせごとに、最良の部分経路スコアを持つものからなるリストが求められる。次

に、その先頭要素が変数parseに、残りの要素が変数parsesに格納される。ここで、関数firstはリストの先頭要素を返し、関数restはリストから先頭要素を除いたリストを返すものとする。

```
05 parses:=get-parses(i, parse-table);
   parse:=first(parses);
   parses:=rest(parses);
```

ステップ4) 最後の2つの単語の品詞の異なる組み合わせ毎に、最適部分経路スコアを持つパース構造のリストの終わりに達したかどうかを判断する。もし、そうなら、ステップ12において文字位置のインデックスをインクリメントする。そうでなければ、ステップ5以降の処理を各パース構造に対して行う。

【0076】parse==nil?

15 ステップ5) 関数get-leftmost-substringsを用いて文字位置iから始まる部分列と一致する辞書中の単語のリストが求められ、変数wordsに代入される。次にその先頭要素が変数wordに、残りの要素が変数wordsに代入される。

```
20 words:=get-leftmost-substrings(string, i);
   word:=first(words);
   words:=rest(words);
```

ステップ6) 文字位置iから始まる部分列と一致するすべての単語が調べられたかどうかを検査する。もし全て調べたのであれば、ステップ11においてparsesの先頭要素をparseに代入し、parsesの先頭要素を取り除く。そうでなければ、ステップ7以降の処理を行う。

【0077】word==nil?

ステップ7) 変数parse, nth-order-stateの値と変数wordのposスロットの値から、品詞の三つ組を表すリストを作成し、変数pos-ngramに代入する。ここで、関数listは引数を要素とするリストを返し、関数appendは引数のリストに含まれる要素を連結したリストを返すものとする。

```
35 【0078】pos-ngram:=append(parse. nth-order-state, list(word. pos));
```

ステップ8) 変数pos-ngramに代入された品詞三つ組の確率が調べられる。もし、品詞三つ組確率が0ならば、ステップ10において、次の要素を変数wordへ代入する。もし、品詞三つ組確率が0でなければ、以下の処理を行う。ここで、関数transprobは引数の品詞三つ組の確率を返すものとする。

【0079】transprob(pos-ngram)=0

ステップ9) まず、新しいパース構造をつくり変数new-parseに代入する。関数make-parseは新しいパース構造を返すものとする。この新しいparse構造の開始位置(new-parse.start)は、i、終了位置(new-parse.end)は、iと変数wordに格納された単語の表記の長さの和、品詞(new-parse.pos)は、変数wordに格納された単語の品詞である。また、結合状態(new-parse. nth-order-sta

te) は、品詞三つ組の先頭要素を取り除いたものであり、最適経路スコア(new-parse.prob-so-far) は、直前の単語の最適経路スコア(parse.prob-so-far) と品詞三つ組確率(transprob(pos-ngram)) と単語出力確率(word.prob) の積である。次に、この新しいparse 構造を、関数register-parseを用いてパーステーブル(parse-table) に登録し、ステップ10に移行する。

```
new-parse:=make-parse();
```

```
new-parse.start:=i;
```

```
new-parse.end:=i+length(word.form);
```

```
new-parse.pos:=word.pos;
```

```
new-parse.nth-order-state:=rest(pos-ngram);
```

```
new-parse.prob-so-far:=parse.prob-so-far*transprob(pos-ngram)*word.prob;
```

```
register-parse(new-parse,parse-table);
```

ステップ10) リストwords の先頭要素をwordへ代入し、words の残りの要素をwords へ代入する。

```
word:=first(words);
```

```
words:=rest(words);
```

ステップ11) リストpasesの先頭要素をparse へ代入し、pasesの残りの要素をpasesへ代入する。

```
parse:=first(words);
```

```
pases:=rest(pases);
```

$$f(n) = g(n) + h(n)$$

本発明は、後向き探索では、関数gとして文末から現在の単語(パース構造)に至るまでの品詞三つ組確率と品詞別単語出力確率の積の対数の絶対値を用いる。また、関数hとしては、文頭から現在に至るまでの品詞三つ組確率と品詞別単語出力確率の積の最大値の対数の絶対値を用いる。

【0084】この後向きA*探索のために図6に示すようなスロット(フィールド)を持つパス(path)という構造体を定義する。構造体pathはA*探索におけるグラフのノードに相当し、現在の単語(parse構造)、これまでの経路及びコストに関する情報を保持する。Path.parseは、parse構造を格納する。Path.previousは直前のpath構造へのポインタである。Path.cost-so-farは、初期状態からのコストである。Path.total-costは、初期状態から最終状態までのコストの推定値である。

【0085】A*探索では、“open”と“close”という二つのリストを用いる。リスト“open”は、既に生成され、ヒューリスティック関数が適用されているが、まだ展開されて(調べられて)いないノード(path構造)の集合である。このリストはヒューリスティック関数の値に基づく優先度付きキュー(priority queue)になっている。リスト“close”は、既に展開された(調べられた)ノードの集合である。

【0086】A*探索では、目標状態に対応するノードを生成するまで、各ステップの一つのノードを展開する。各ステップでは、既に生成されているが、まだ、展

ステップ12) 変数iの値をインクリメントする。

```
【0080】 i=i+1
```

(4) A*アルゴリズムを用いた後向き探索

後向きA*探索部2が最も尤もらしい順にひとつずつ形態素解析候補を求める手順を示す。最初に、A*探索の概要、後向きA*探索部で用いられるデータ構造、及び、補助的な関数について説明し、続いて、具体的な処理の流れを説明する。

【0081】本発明の後向きA*探索では、単語と品詞の組であるパース構造を、A*アルゴリズムにおけるグラフのノードと考える。そして、コストとしては、確率の対数の絶対値を用いる。これにより、確率最大の解は、コスト最小の解に対応し、確率の積はコストの和に対応する。

【0082】A*探索では、ヒューリスティック関数f(n)を考える。ヒューリスティック関数f(n)は、現在のノードnを生成した経路に沿って、初期状態から最終状態へ至るまでのコストの推定値を与える。初期状態から現在のノードへ至るまでのコストを与える関数をg(n)、現在のノードから最終状態へ至るまでのコストの推定値を与える関数をh(n)とすると、ヒューリスティック関数f(n)は次式により与えられる。

```
【0083】
```

(9)

開されていない最も有望なノードを展開する。即ち、選ばれたノードの後続のノード(successors)を生成し、ヒューリスティック関数を適用し、既に生成されていないかを検査した後にリスト“open”に加える。この検査によって、各ノードはグラフの中に一回だけ現れることが保証される。また、二つ以上の経路が同じノードに生成する時は、スコアの良方を記録する。

【0087】リスト“open”とリスト“close”を操作する補助関数として、“find-path, insert-and-sort-path, delete-path”を定義する。関数“find-path”はパース構造とリスト(open又はclose)を引数とし、引数のパース構造と、次に定義する意味で「等しい」パース(parse)構造を持つパス(path)構造がリスト中に存在すれば、それを返す。ここで、二つのパース構造の開始位置と終了位置と結合状態が等しい場合、この二つのパース構造は「等しい」と定義する。

【0088】関数“insert-and-sort-path”はパス(path)構造とリストを引数とし、リストにパスを挿入し、初期状態から最終状態までのコストの推定値“path.total-cost”の小さい順にソートする。関数“delete-path”は、パス構造とリストを引数とし、引数のパス構造をリストから削除する。

【0089】また、後向き探索において、ノード(パス構造)の後続ノードを生成するための補助関数として、“immediate-left-pases”を定義する。関数“immediate-left-pases”は、パース構造とパーステーブル(

最適経路スコアテーブル)を引数とし、引数のパース構造の左側(文頭側)に接続しうるすべてのパース構造のリストを返す。具体的には、次に三つの条件を満たすパース構造をパーステーブルの中から検索する。

(1) 引数のパース構造の開始位置がこのパース構造の終了位置である:

(2) 引数のパース構造の結合状態のリストが最後の要素を無視したものと等しい:

(3) このパース構造から引数のパース構造への前向きの状態遷移を表す品詞三つ組の確率が0ではない:さらに、後向き探索において、初期状態から現在のノードへ至るまでのコストを計算するために、ノード間の遷移コストを計算する関数として“transition-and-word-cost”を定義する。関数“transition-and-word-cost”は、隣接する二つのパース構造を引数(第1引数が文末側の単語で、第2引数が文頭側の単語)とし、文末側の単語の品詞別単語出力確率の対数と品詞三つ組確率の対数の和の絶対値を返す。

【0090】また、現在のノードから最終状態へ至るまでのコストを得るために、関数“cost-from-beginning-of-sentence”を定義する。関数“cost-from-beginning-of-sentence”は、文頭から引数のパース構造に至るまでのコストを返す。このコストは、既に前向き探索で既に計算されており、“parse. prob-so-far”に格納されている。

【0091】図9は、本発明の一実施例のA*探索部の動作を説明するためのフローチャートである。以下では、同図に従って、後向きA*探索部の動作を説明する。

【0092】ステップ101) 文末において、最適経路コストをもつパース構造をパーススロットに持つパース構造からなるリストを“open”に代入する。ここで、関数“backward-search-initial-paths”は、このような初期化を行う関数とする。また、“closed”には、空リストを代入する。

```
open:=backward-search-initial-paths();
closed:=nil;
```

ステップ102) “open”が空リストかどうかを調べる。もし空リストであれば、解がみつからなかったので探索が失敗したことを通知して探索を終了する。そうでなければ、ステップ103以降の処理を行う。

【0093】open==nil ?

ステップ103) “open”の先頭要素を変数“bestpath”に代入する。

【0094】bestpath:=first(open)

ステップ104) 探索が文頭に達したかを調べる。ここで、関数“is-beginning-of-sentence”は、この条件を調べる関数である。もし、探索が文頭に達していれば、“bestpath”が最適解であり、探索が成功したことを通知して探索を終了する。そうでなければ、ステップ105

以降の処理を行う。また、探索は成功したが、さらにその次に最も尤もらしい解を求めたい場合にも、ステップ105以降の処理を行う。

【0095】

05 is-beginning-of-sentence(bestpath-parse) ?

ステップ105) まず、変数“open”へ“open”から先頭要素を取り除いたリストを代入する。次に、関数“insert-and-sort-path”を用いて、“bestpath”をリスト“closed”へ挿入し、“closed”の要素を初期状態から最終状態までのコストの推定値に順にソートする。

【0096】open:=rest(open);

insert-and-sort-path(bestpath, closed)

ステップ106) 関数“immediate-left-parses”を用いて、“bestpath”のパーススロットに格納されたパース構造の左側に接続するすべてのパース構造をパーステーブルから取り出したリストを作り、このリストを“successorparses”に代入する。続いて、successor-parseの先頭要素を“successorparse”に代入し、残りを“successorparses”に代入する。

20 【0097】successorparses:=immediate-left-parses(bestpath, parse, parse-table);

successorparse:=first(successorparses);

successorparses:=rest(successorparses);

ステップ107) “successorparses”のすべての要素を処理したかを検査する。もし、そうであれば、ステップ102に移行し、そうでなければステップ108の処理を行う。

【0098】successorparse==nil ?

ステップ108) 新しいパース構造を割り当て、これを変数“newpath”に代入する。関数“make-path”は新しいパース構造のための記憶領域を割り当てる関数である。そして、“newpath.parse”には、“successorparse”を代入し、“newpath.previous”には“bestpath”を代入する。“newpath.cost-so-far”には、“bestpath.parse”へ至るまでのコスト“bestpath.cost-so-far”と、“bestpath.parse”から“newpath.parse”への遷移のコストの和が代入される。この遷移コストは、関数“transition-and-word-cost”を用いて計算される。

“newpath.total-cost”には、文末から“newpath.parse”までのコスト(newpath.cost-so-far)と、文頭からnewpath.parseまでのコストの和が代入される。文頭からのコストは関数“cost-from-beginning-of-sentence”を用いて計算される。

newpath:=make-path();

45 newpath.parse:=successorparse;

newpath.previous:=bestpath;

newpath.cost-so-far:=bestpath.cost-so-far

+transition-and-word-cost(bestpath.parse, newpath.parse);

50 newpath.total-cost-so-far

+cost-from-beginning-of-sentence(newpath.parse);
 ステップ109) 関数“find-path”を用いて、successorparse をパーススロットにもつパス構造が“open”に含まれているかを検査する。もし含まれていなければ、ステップ113に移行する。含まれていれば、ステップ110以降の処理を行う。

【0099】find-path(successorparse, open)==nil ?
 ステップ110) “successorparse”をパーススロットに持つ“open”の中のパス構造を変数“oldpath”に代入する。

【0100】
 oldpath:=find-path(successorparse, open);
 ステップ111) “newpath.total-cost”と“oldpath.total-cost”を比較する。もし、“newpath.total-cost”の方が大きければ、何もせずにステップ118に移行する。もし、“newpath.total-cost”の方が小さければ、ステップ112に移行する。

【0101】
 newpath.total-cost<oldpath.total-cost ?
 ステップ112) 関数“delete-path”を用いて“open”から“oldpath”を削除し、関数“insert-and-sort-path”を用いて、“newpath”を“open”へ挿入した後にソートする。そしてステップ118に移行する。

【0102】delete-path(oldpath, open);
 insert-and-sort-path(newpath, open);
 ステップ113) 関数find-pathを用いてsuccessorparseをパーススロットに持つパス構造が“closed”に含まれているかを検査する。もし含まれていれば、ステップ117に移行する。含まれていなければステップ114以降の処理を行う。

【0103】
 find-path(successorparse, closed)==nil ?
 ステップ114) “successorparse”をパーススロットに持つ“closed”の中のパス構造を変数“oldpath”に代入する。

【0104】
 oldpath:=find-path(successorparse, closed);
 ステップ115) “newpath.total-cost”と“oldpath.total-cost”を比較する。もし、“newpath.total-cost”の方が大きければ、何もせずにステップ118に移行する。

【0105】newpath.total-cost<oldpath.total-cost
 ステップ116) ステップ115においても、“newpath.total-cost”の方が小さければ、関数“delete-path

h”を用いて、“closed”から“oldpath”を削除、関数“insert-and-sort-path”を用いて“newpath”を“closed”へ挿入した後にソートする。そしてステップ118に移行する。

05 delete-path(oldpath, closed);
 insert-and-sort-path(newpath, open);
 ステップ117) 関数“insert-and-sort-path”を用いて“newpath”を“open”へ挿入した後にソートする。そしてステップ118に移行する。

10 【0106】insert-and-sort-path(newpath, open);
 ステップ118) “successorparse”の先頭要素を“successorparse”へ代入し、“successorparse”の残りの要素を“successorparse”へ代入する。そして、ステップ107に移行する。

15 successorparse:=first(successorparse);
 successorparse:=rest(successorparse);
 図3の後向きA*探索部2により得られた形態素解析候補の例において、図10の各文に対して上記3個の形態素解析候補が示されている。各形態素解析候補には、確率の対数が示されており、この値が大きいほど尤もらしい。

【0107】図10において、『会議に申し込みたいのですが。』が入力文として与えられたときに、第1候補の確率が自然対数をとったものが、-25.82409083887518であり、この形態素解析候補は、「会議」「に」「申し込み」「たい」「の」「です」「が」「。」という単語系列と、「普通名詞」「格助詞」「本助詞・連用・五段」「助動詞・連体」「準体助詞」「助動詞・終止」「接続助詞」「記号」という品詞系列から構成されることを表している。第2候補及び第3候補についても同様である。確率は、自然対数をとっているため、負の値となるが、値が大きいほど(絶対値が小さいほど)尤もらしい候補となる。

【0108】(5) 単語モデルの推定
 35 単語モデル推定部7が、品詞タグ付きコーパスから品詞別文字三つ組確率から構成される単語モデルを推定する手順を示す。

【0109】本発明では、品詞別の文字三つ組確率を用いて、単語モデルを作成する。単語wを構成する文字列を $C=c_1c_2\cdots c_n$ とする。品詞tが与えられた時の単語の出力確率を次式で近似させる。ここで“#”は単語境界を示す。

【0110】
 【数5】

$$P(w|t) = P_1(C) = P_1(c_1 | \#, \#) P_1(c_2 | \#, c_1)$$

$$\prod_{i=1}^n P_1(c_i | c_{i-2}, c_{i-1}) P_1(\# | c_{n-1}, c_n)$$

(10)

【0111】品詞別文字三つ組確率は、品詞タグ付きコーパスにおいて、品詞が t である全ての単語に現れた文字三つ組と文字二つ組の相対頻度から次式により推定できる。ここで、 $N_1(c_{i-2}, c_{i-1}, c_i)$ は、品詞タグ t が与えられた単語に現れた文字三つ組 c_{i-2}, c_{i-1}, c_i の総数である。この相対頻度は、単語の出現頻度を反映することに注意して欲しい。

$$P_1(c_i | c_{i-2}, c_{i-1}) = f_1(c_i | c_{i-2}, c_{i-1})$$

$$= \frac{N_1(c_{i-2}, c_{i-1}, c_i)}{N_1(c_{i-2}, c_{i-1})}$$

(11)

【0113】一般に、この文字三つ組確率を推定するのに十分な量の訓練データを用意することは難しい。そこで

$$P(c_i | c_{i-2}, c_{i-1}) = q_2 f(c_i | c_{i-2}, c_{i-1}) + q_1 f(c_i | c_{i-1}) + q_0 f(c_i) + q_0 V \quad (12)$$

各品詞別単語モデルは、ある文字列がその品詞の単語である場合の単語出力確率を求める。従って、単語仮説の文字列表記と品詞が与えられれば、この単語モデルから品詞別単語出力確率を求めることができる。なお、品詞別単語モデル（文字三つ組モデル）は、名詞や動詞等の全てのオープンカテゴリ、即ち、その品詞に属する単語の数が有限ではないと考えられる品詞に対して作成する。

【0114】(6) 単語仮説の生成

単語仮説生成部 9 が、単語モデルを用いて入力文に対して単語仮説を生成する手順を説明する。

【0115】単語仮説生成部 9 は、未知語を扱うために、前向き探索において、入力文のある文字位置から始まる任意の長さの部分文字列と全てのオープンカテゴリ（品詞）の組に対して、単語モデルを用いて品詞別単語出力確率を計算し、最も尤もらしい順に予め決められた個数の単語仮説を生成する。

【0116】図 11 は、単語仮説生成部の動作を説明するためのフローチャートである。以下では、同図に従って、単語仮説生成部 9 の動作を説明する。

【0117】単語仮説生成部 9 は、入力文字列 “string”、及び単語仮説を生成すべき文字位置 i を引数として与えられる。また、定数 “delimiters” には、句点（。）や読点（、）等の単語の一部とはなり得ない区切り記号のリストが格納され、定数 “open-categories” には、名詞や動詞などのその品詞に属する単語の数が有限ではない品詞（オープンカテゴリ）のリストが格納さ

れている。この相対頻度は、単語の出現頻度を反映することに注意して欲しい。

【0112】

【数 6】

で、品詞三つ組の場合と同様に、次式に示すように内挿推定法を用いて文字三つ組確率を平滑化する。

【0118】ステップ 201）単語仮説のリストを格納する変数 “word-hypos” に初期値として “nil” を代入する。

【0119】word-hypos:=nil;

ステップ 202）単語仮説の長さを格納する変数 j に初期値として 0 を代入する。

【0120】 $j := 0$

ステップ 203）単語仮説の長さ i が単語仮説を生成すべき文字位置 i から文末までの文字列の長さより大きくなったかどうかを調べる。もし、そうならば、ステップ 209 に移行する。そうでなければ、ステップ 204 以降の処理を行う。

【0121】 $j > \text{length}(\text{string}) - i$

ステップ 204）単語仮説となる文字列の最後の文字が区切り記号であるかどうかを調べる。もし、そうならば、ステップ 209 に移行する。区切り記号でなければ、ステップ 205 以降の処理を行う。

【0122】member(char(string, $i+j-1$), delimiters)

ステップ 205）入力文字列 “string” の位置 i から始まる長さ j の部分文字列を変数 form へ代入する。ここで、関数 substr は、部分列を返す関数とする。

【0123】form:=substr(string, i, j);

ステップ 206）オープンカテゴリリスト “open-categories” の先頭要素を変数 pos に代入し、残りを変数 “poses” に代入する。

【0124】poses:=open-categories;


```
pos:=first(poses);
poses:=rest(poses);
```

ステップ 207) オープンカテゴリリストの終わりに達したかどうかを調べる。もしそうであれば、ステップ 203 に移行する。そうでなければステップ 208 以降の処理を行う。

【0125】 pos==nil ?

ステップ 208) まず、品詞別の文字三つ組確率を用いて、表記が“form”で品詞が“pos”である単語の出現確率を求め、変数“prob”に代入する。関数“wordprob-with-word-model”は、表記と品詞を引数とし、品詞別文字三つ組確率から構成される単語モデルを用いて求めた品詞別単語出力確率を返す。次に、表記と品詞と確率から構成されるリストを変数“word-hypo”に代入し、これを単語仮説リスト“word-hypos”の先頭に加え、ステップ 207 に移行する。

```
prob:=wordprob-with-word-model(from, pos);
word-hypo:=list(form, pos, prob);
push(word-hypo, word-hypos);
```

ステップ 209) 単語仮説リスト“word-hypos”を品詞別単語出力確率の大きい順にソートする。

【0126】 sort-by-prob(word-hypos);

図 11 は、引数として与えられた文字列の左端から始まる部分列に対して単語仮説を生成した例を示す。

【0127】 同図において、入力文に未知語がある場合でも出力は、上記の図 9 のような形式になる。図 11 の例では『転送』という単語がシステムには登録されていないが、単語モデルに基づいて生成された単語仮説の中には、『転送』が「サ変名詞」であるような正解も含まれている。最終的な形態素解析結果は本発明の前向き探索、後向き探索を用いて未知語の前後の部分との関係を考慮した上で決定するため、単語仮説生成の段階で正解が一候補である必要はない。

【0128】 単語が全て辞書に登録されている場合に最尤候補（最も確率が高い形態素解析候補）を求めるだけであれば、後向き A* 探索は必要ないが、後向き A* 探索（グラフ探索）を用いることにより、辞書に未登録の単語も形態素解析の対象とすることが可能となる。

【0129】 なお、上記の実施例では、統計的言語モデルとしてトライポス (tri-POS) モデルを使用した。この例に限定されることなく、トライクラス (triclass)、トライ・タグ (tri-tag)、トライ・グラム (tri-gram) 等を用いてもよい。

【0130】 なお、本発明は、上記実施例に限定されることなく、特許請求の範囲内で種々変更及び応用が可能である。

【0131】

【発明の効果】 上述のように、本発明によれば、品詞タグ付きコーパスから推定した品詞三つ組確率と品詞別単語出力確率から構成される品詞付きモデル、品詞タグ付

コーパスから推定した品詞別文字三つ組確率から構成される単語モデル、単語モデルに基づく単語仮説生成、動的計画法を用いた前向き探索、及び A* アルゴリズムを用いた後向き探索により品詞列と単語列の同時確率を最大化する入力文の形態素解析候補、即ち、単語分割と品詞付与の組を入力文が未知語を含む場合でも最も尤もらしい候補から順に任意の数だけ出力できるので、高い精度を持ち、頑強でかつ柔軟なインタフェースを持つ日本語形態素解析装置を実現できる。

10 【図面の簡単な説明】

【図 1】 本発明の原理構成図である。

【図 2】 本発明の原理を説明するためのフローチャートである。

15 【図 3】 本発明の一実施例の日本語形態素解析装置のモジュール図である。

【図 4】 タグ付コーパスの例を示す図である。

【図 5】 品詞三つ組確率の例を示す図である。

【図 6】 品詞別単語出力確率を示す図である。

20 【図 7】 本発明の一実施例の N-best アルゴリズムのためのデータ構造を示す図である。

【図 8】 本発明の一実施例の前向き DP 探索を示すフローチャートである。

【図 9】 本発明の一実施例の後向き A* 探索を示すフローチャートである。

25 【図 10】 本発明の一実施例の形態素解析候補の例を示す図である。

【図 11】 本発明の一実施例の単語仮説を生成する動作のフローチャートである。

【図 12】 本発明の一実施例の単語仮説の例を示す図である。

【符号の説明】

1 前向き DP 探索部

2 後向き A* 探索部

3 品詞タグ付コーパス

35 4 品詞付けモデル推定部

5 品詞三つ組確率テーブル

6 品詞別単語出力確率テーブル

7 単語モデル推定部

8 品詞別文字三つ組確率テーブル

40 9 単語仮説生成部

10 最適経路スコアテーブル

100 品詞付けモデル

110 品詞付モデル推定手段

120 品詞三つ組確率テーブル

45 130 品詞別単語出力確率テーブル

140 前向き探索部

150 後向き探索部

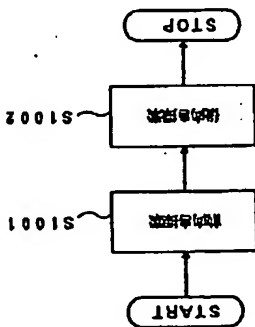
170 最適経路スコアテーブル

180 品詞タグ付きコーパス

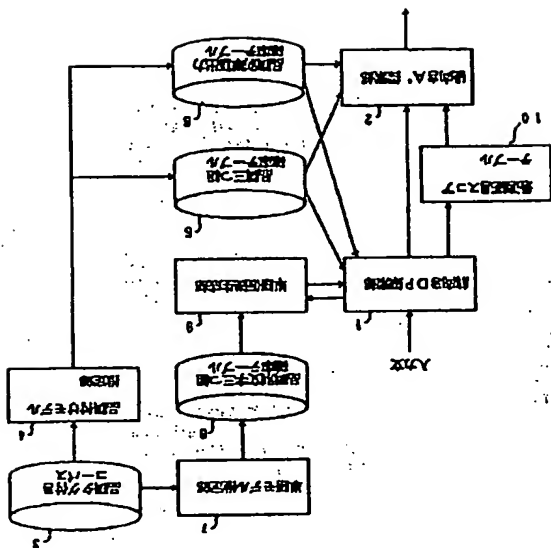
50 190 単語モデル

【図2】

本書の目的は、このように、

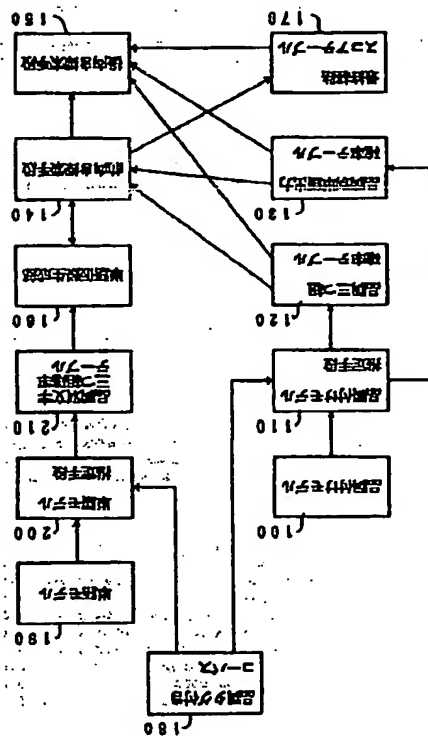


【図 3】

[illegible]

【 1 图】

● 2010年10月1日



【9区】

图 6-1-2 本例 4 个开闭状态

[illegible]

【図4】

品名タグ付きコーバスの明を示す図

3001-100
もしし／所動詞 へ／12号 通ひ／所動詞 通ひ／所動詞 通ひ／所動詞 通ひ だす／所動詞・終止
か／所動詞 ？／12号

3001-200
はい／所動詞 。／12号 せう／所動詞 だす／所動詞・終止 。／12号

3001-300
全編／所動詞 に／所動詞 申ひ／所動詞・通用・五段 たい／所動詞・連体
の／所動詞 だす／所動詞・終止 か／所動詞 。／12号

3001-400
はい／所動詞 へ／12号 書け／所動詞 用紙／所動詞 は／所動詞 閉じ／所動詞
お／所動詞 持ち／所動詞・通用・五段 だし／所動詞・未然 う／所動詞・終止
か／所動詞 ？／12号

3001-500
はい／所動詞 。／12号 まだ／所動詞 だす／所動詞・終止 。／12号

【図5】

品目三つ相対率を示す図

[illegible]

【圖 7】

本発明の一実施例のN-bestアルゴリズムのための
データ構造を示す図

parse 構造	
start end pos nth-order-state prob-so-far	入力文字列における単語の開始位置 入力文字列における単語の終了位置 単語の品詞 表現の二つの単語の品詞のリスト 文節からの最長部分結合スコア
word 構造	
form pos prob	単語の表記 単語の品詞 品詞内の単語の出力確率
path 構造	
parse previous cost-so-far total-cost	parse 構造 直前の path 構造へのポイント 初項状態からのコスト 初項状態から最終状態までのコスト

【图 10】

本邦四月の一実施例の形態を併列示す図

> (memberlist "全国に申し込みたいのですが。")
 -1. 42409032387518
 全国／普通会員 に入会費 申し込み／本会費／通用・五段 たし／加算費・連体の／年次会費
 です／加算費・禁止 あり／加算費。／代号
 -2. 253535232513255
 全国／普通会員 に入会費 申し込み／普通会員 たし／加算費・連体の／年次会費
 です／加算費・禁止 あり／加算費。／代号
 -3. 4232072321427
 全国／固有会員 に入会費 申し込み／本会費／通用・五段 たし／加算費・連体の／年次会費
 です／加算費・禁止 あり／加算費。／代号

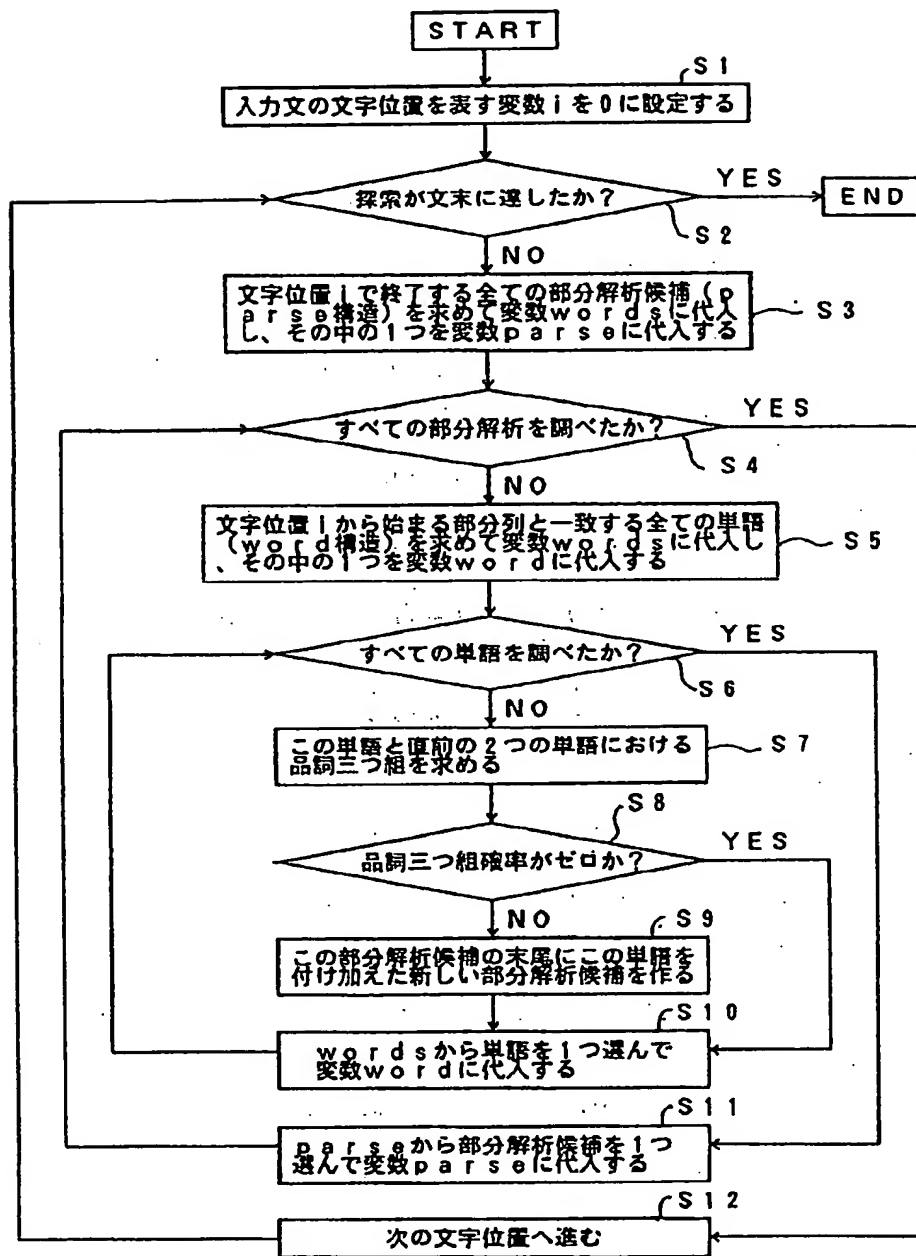
【圖 12】

本発明の一実施例の半導体装置の例を示す図

```
> (get-leftmost-substrings-all chr:model "転送して667下8い")
((転送 サ変名) 2.51945769735981E-7)
(転送 国名) 2.344921507010935E-6)
(転送 書道名) 7.0243990747133745E-9)
(転送 姓) 2.3758509760988E-6)
(転送 サ変名) 5.70817493025141E-12)
(転送 国名) 6.928642348107103E-14)
(転送 サ変名) 1.824010652511430E-14)
(転送 書道名) 1.1267031867281895E-14)
(転送 姓) 6.8994934913207E-10)
(転送 サ変名) 1.94063591182284E-18)
> (get-leftmost-substrings-all chr:model "転送して667下8い")
((転送 サ変名) 2.51945769735981E-7)
(転送 国名) 2.344921507010935E-6)
(転送 書道名) 7.0243990747133745E-9)
(転送 姓) 2.3758509760988E-6)
(転送 サ変名) 5.70817493025141E-10)
(転送 書道名) 4.736222004870369E-13)
(転送 国名) 6.928642348107103E-14)
(転送 サ変名) 7.236813394426545E-14)
(転送 姓) 6.8994934913207E-10)
(転送 書道名) 2.403023305251351E-17)
```

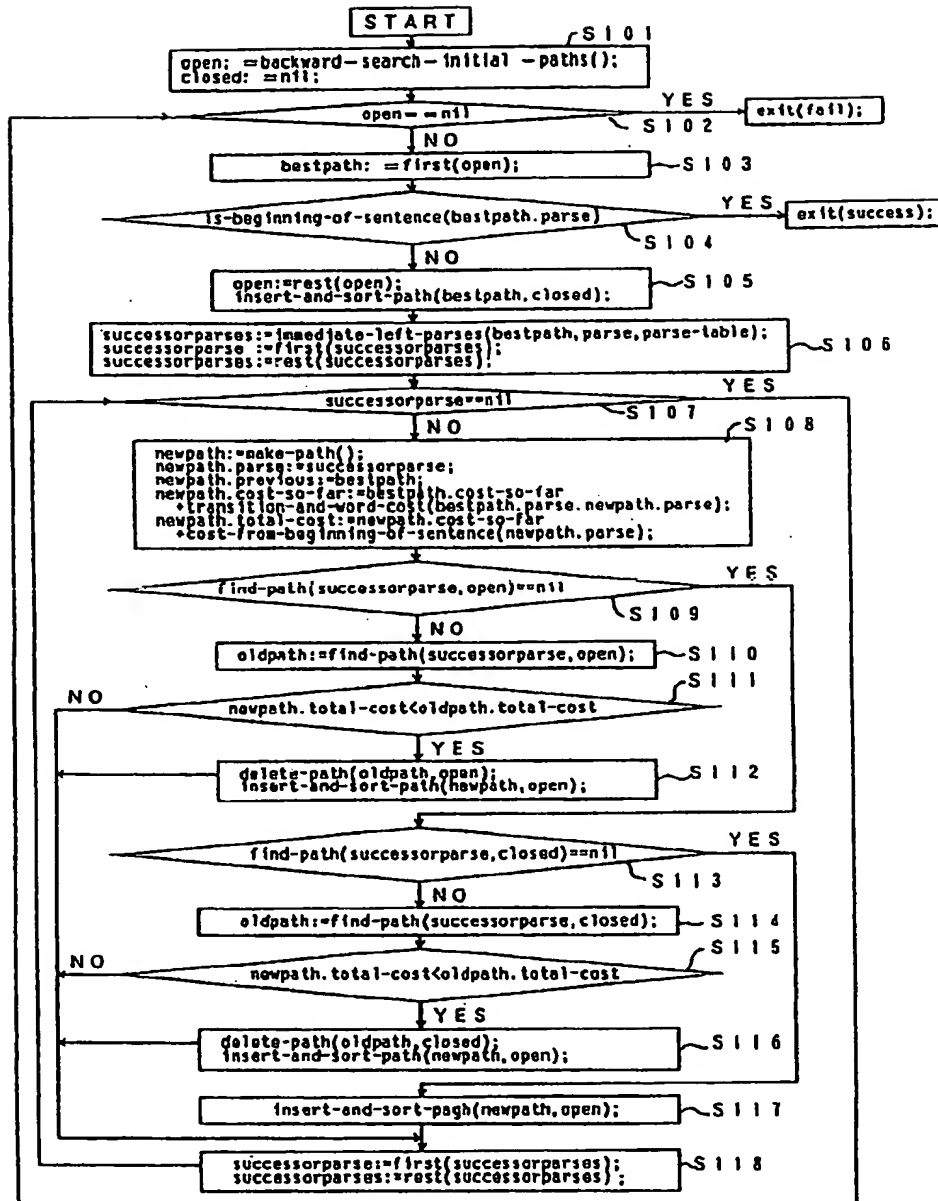
【図 8】

本発明の一実施例の前向き DP 探索を示すフローチャート



【図9】

本発明の一実施例の逆向きA*探索を示すフローチャート



【図 11】

本発明の単語仮説を生成する動作のフローチャート

